

DS2000

4/11 - Tues.

## Admin

- HW8 should be returned Thurs
- no class on 4/18 :)
- this Fri (4/14) - wrap up day

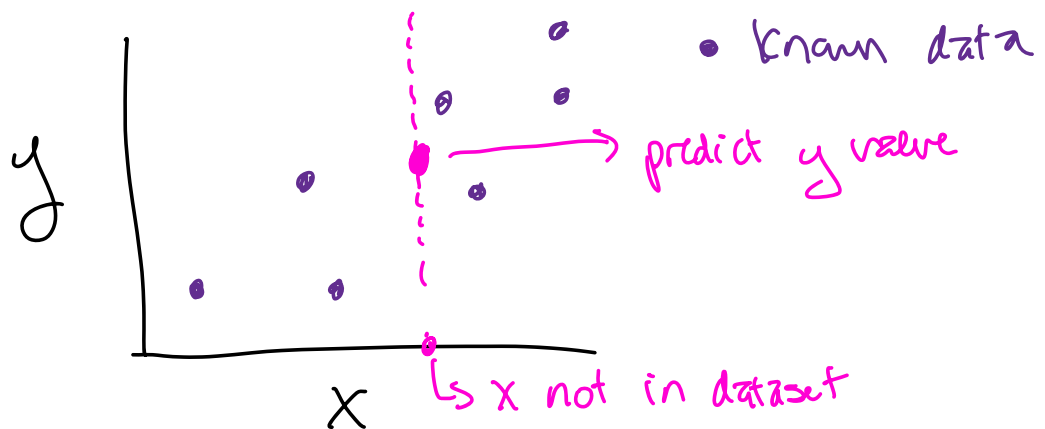
[bit.ly/ds2000-wrap](http://bit.ly/ds2000-wrap)

## Agenda

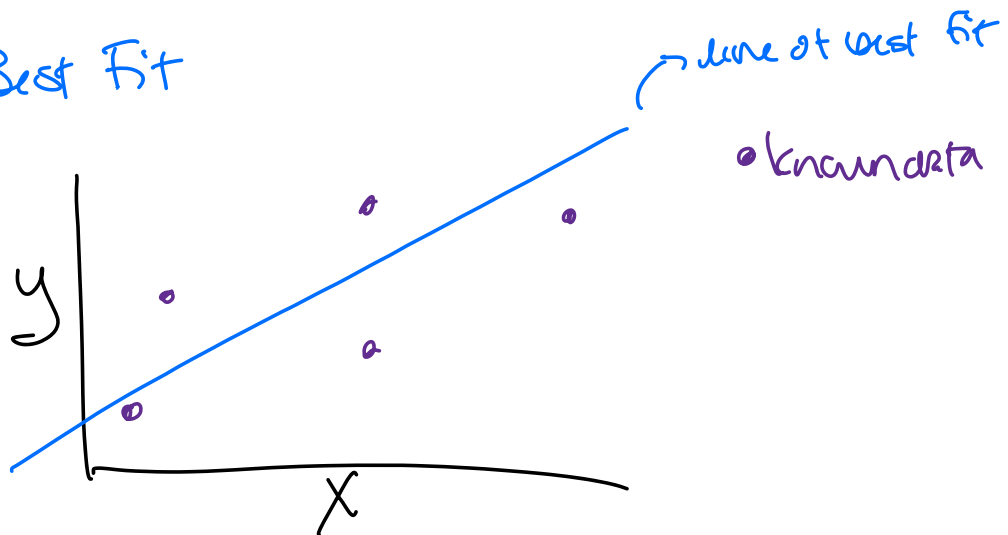
1. Linear regression
2. Pandas, seaborn, scipy
3. Python

# 1. Linear Regression

- relationship between 2 variables
  - Used as prediction model
  - $x$ : independent variable
  - $y$ : dependent variable
- } given  $x$  not in dataset, what would  $y$  be?



## Line of Best Fit

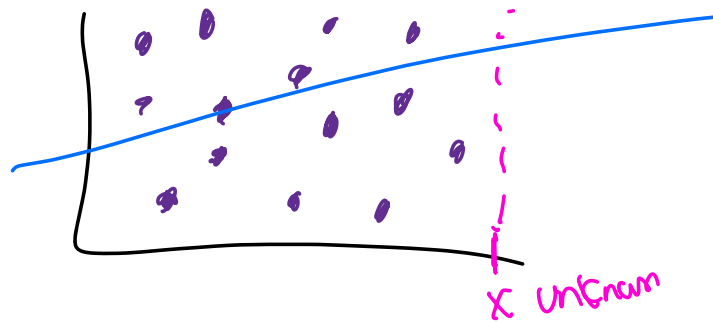


want to predict  $y$ , given  $x$

$$y = \underset{\uparrow}{m}x + \underset{\uparrow}{b}$$

↳ slope ↳ intercept

Python will give a line  
of best fit for any data



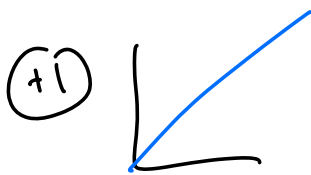
But : First step!

- What is correlation between  $x$  and  $y$ ?
- If highly correlated, using the line of best fit makes sense
- default in Python: pearson

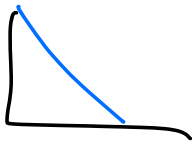
↓  
-1      to      +1

close to zero: bad ;)

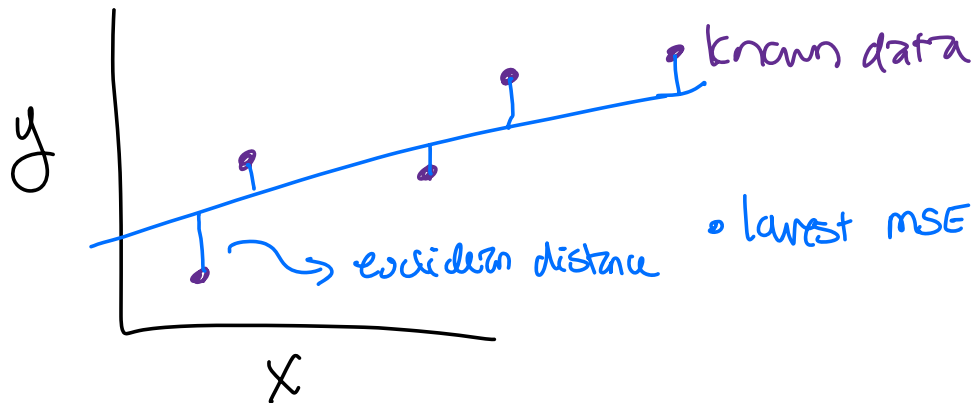
close to +1, -1: good ;)



①



Line of Best Fit : mean squared error (MSE)



## 2. Pandas/seaborn/scipy

Today's data:

- Nicolas (age movies, precipitation in Boston)
- JP property values (time series)

### Python we need

Pandas — file, correlation

```
import pandas as pd
```

```
df = pd.read_csv(~~)
```

```
df["Precipitation"] ~~~> 2 column
```

```
df.corr() ~~~~~> correlation
```

Seaborn — plot the line

```
import seaborn as sns
```

```
sns.regplot(~~) ~~~> draws  
Scatter + line
```

Scipy — get the slope, intercept

```
from scipy import stats
```

```
stats stats.linregress() ~~~> get  
slope/int.
```